**Original Article**

# Details of content validity and objectifying it in instrument development

Vahid Zamanzadeh[1], Maryam Rassouli[2], Abbas Abbaszadeh[3], Hamid Alavi Majd[4], Alireza Nikanfar[5], Akram Ghahramanian[1*]

[1] Department of Medical Surgical, School of Nursing and Midwifery, Tabriz University of Medical Sciences, Tabriz, Iran

[2] Department of Pediatrics, School of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran

[3] Department of Medical Surgical, School of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran

[4] Department of Biostatistics, School of Para Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran

[5] Department of Internal Diseases, School of Medicine AND Hematology & Oncology Research Center, Tabriz University of Medical Sciences, Tabriz, Iran

## ARTICLE INFO

## ABSTRACT

**Background & Aim:** Researchers in the nursing science study complex constructs for which valid and reliable instruments are needed. When an instrument is created, psychometric testing is required, and the first-step is to study the content validity of the instrument. This article focuses on the process used to assess the content validity.

**Methods & Materials:** This article examines the definition, importance, conceptual basis, and functional nature of content validity in instrument development. The conditional and dynamic nature of content validity is discussed, and multiple elements of content validity along with quantitative and qualitative methods of content validation are reviewed.

**Results:** In content validity process, content representativeness or content relevance of the items of an instrument is determined by the application of a two-stage (development and judgment) process. In this review, we demonstrate how to conduct content validity process, to collect specific data for items generation and calculation of content validity ratio, content validity index, modified Kappa coefficient, and to guide for interpreting these indices. Face validity through suggestions of expert panel and item impact scores is also discussed in paper.

**Conclusion:** Understanding content validity is important for nursing researchers because they should realize if the instruments they use for their studies are suitable for the construct, population under study, and sociocultural background in which the study is carried out, or there is a need for new or modified instruments.

## Introduction

In most studies, researchers study complex constructs for which valid and reliable instruments are needed (1). Validity, which is defined as the ability of an instrument to measure the properties of the construct under study (2), is a vital factor in selecting or applying an instrument. It is determined as its three common forms including content, construct, and criterion-related validity (3). Since content validity is a prerequisite for other validity, it should receive the highest priority during instrument development. Validity is not the property of an instrument, but the property of the scores achieved by an instrument used for a specific purpose on a special group of respondents. Therefore, validity evidence should be obtained on each study for which an instrument is used (4).

Content validity, also known as definition validity and logical validity (5), can be defined as

* Corresponding Author: Akram Ghahramanian, Postal Address: School of Nursing and Midwifery, Tabriz University of Medical Sciences, Sought Shariati Street, Tabriz, Iran.
Email: ghahramaniana@gmail.com

the ability of the selected items to reflect the features of the construct in the measure. This type of validity addresses the degree to which items of an instrument sufficiently represent the content domain. It also answers the question that to what extent the selected sample in an instrument or instrument items is a comprehensive sample of the content (1, 6-8). This type validity provides the preliminary evidence on construct validity of an instrument (9). In addition, it can provide information on the representativeness and clarity of items and help improve an instrument through achieving recommendations from an expert panel (6, 10). If an instrument lacks content validity, it is impossible to establish reliability for it (11).

Although more resources should be spent for a content validity study initially, it decreases the need for resources in the future reviews of an instrument during psychometric process (1). For this purpose, researchers may gain invaluable information using an expert panel and their feedback. In fact, content validity provides an objective criterion to evaluate each item in an instrument and the entire instrument (3, 12).

Despite the fact that in instrument development, content validity is a critical step (13) and a trigger mechanism to link abstract concepts to visible and measurable indices (7), it is studied superficially and transiently. This problem might be due the fact that the methods used to assess content validity in nursing literature are not referred to profoundly (13), and sufficient details have rarely been provided on content validity process in a single resource (14). It is particularly possible that students do not realize the required complexities in this critical process (13).

On the other hand, a number of experts have questioned historical legitimacy of content validity as a real type of validity (15-17). These challenges about the value and merit of content validity have arisen from lack of distinction between content validity and face validity, unstandardized mechanisms to determine content validity and the previously its un-quantified nature (3). This article aims to discuss the values of content validity by differentiating content validity from face validity, providing a comprehensive process to determine content validity

and introducing possible solutions for quantification of content validity to students and novice researchers in the instrumentation field.

The assessment of content validity begins in the earliest development of an instrument. Researchers can receive invaluable information by conducting a content validity process. Using a panel of experts provides constructive feedback about the quality of the newly developed instrument and objective criteria to evaluate each item (1). We attempt to explain how to achieve the acceptable criteria for content validity of an instrument and to report this process in ourselves study. Content validity is the determination of the content representativeness or content relevance of the items of an instrument using a two-step process (development and judgment). Using a two-stage process is fundamental to determine and quantify content validity in all instrumentation (3). In paper theoretical issues on content validity process are discussed by extensive literature review. In following, this two-step process is discussed. We show the different stages of content validity study in figure 1.

### Stage 1: Instrument development

Instrument development is performed through three steps, including determining content domain, sampling from content (item generation) and instrument construction (11, 15). The first step is determining the content domain of a construct that the instrument is made to measure it. Content domain is the content area related to the variables that being measured (18). It can be identified by literature review on the topic being measured, interviewing with the respondents and focus groups. Through a precise definition on the attributes and characteristics of the desired construct, a clear image of its boundaries, dimensions, and components is obtained. The qualitative research methods can also be applied to determine instrument items (19). The qualitative data collected in the interview with the respondents familiar with concept help enrich and develop what has been identified on the concept, and are considered as an invaluable resource to generate instrument items (20). To determine content domain in emotional instruments and cognitive instruments, we can use literature re-

view and table of specifications, respectively (3). In practice, table of specifications reviews alignment of a set of items (placed in rows) with the concepts forming the construct under study (placed in columns) through collecting quantitative and qualitative evidence from experts and by analyzing data (5). Ridenour et al. also introduced the application of mixed method (deductive-inductive) for conceptualization at the step of the content domain determination and items generation (21). However, generating items requires a preliminary task to determine the content domain for which an instrument is made to measure it. The items can be generated from different resources such as interviewing with experts and respondents to an instrument and literature review (22). In addition, a useful approach would consists of returning to research questions and ensuring that the instrument items are reflect of and relevant to research questions (23).

Instrument construction is the third step in instrument development in which items are refined and organized in a suitable format and sequence so that the finalized items are collected in a usable form (3).
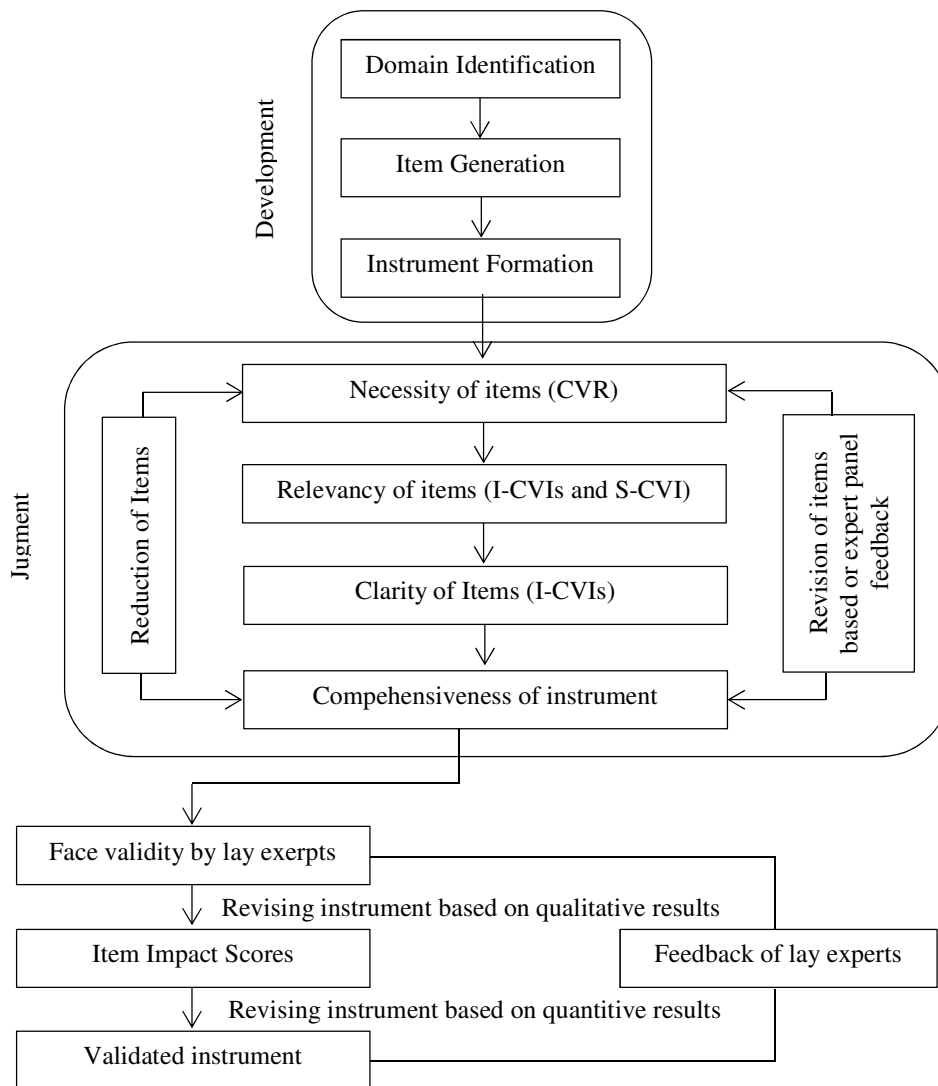


**Figure 1.** Steps of content validity process

### Stage 2: Judgment

This step entails confirmation by a specific number of experts, indicating that instrument items and the entire instrument have content validity. Although investigators often report that the content validity of an instrument is supported by a panel of experts, the characteristics and qualifications of these individuals and the process they are asked to use to assess validity often is not reported (12). In selecting these individuals, was emphasized on the necessity of relevant training, experience, and qualifications of content experts. A history of publications in refereed journals, national presentations, and research on the phenomenon of interest may be used as one criterion in selecting content experts (24). Clinical expertise also may be a criterion used to select panel members. Determining the number of experts has always been partly arbitrary. At least five people are recommended to have sufficient control over chance agreement. The maximum number of judges has not been determined yet; however, it is unlikely that more than ten people are used, but it should be noted as the number of experts increase, the probability of chance agreement decreases. After determining an expert panel, researcher can collect and analyze their quantitative and qualitative viewpoints on the relevancy or representativeness, clarity and comprehensiveness of the items to measure the construct operationally defined by these items to ensure the content validity of the instrument (3, 7, 8).

### Quantification of Content Validity

The content validity of research instruments can be determined using the viewpoints of the panel of experts. This panel consists of content experts and lay experts. Lay experts are the potential research subjects, and content experts are professionals who have research experience or work in the field (25). Using subjects of the target group as expert ensures that the population for whom the instrument is being developed is represented (1).

In qualitative content validity method, content experts and target group's recommendations are adopted on observing grammar, using appropriate and correct words, applying correct and proper order of words in items and appropriate

scoring (26). However, in the quantitative content validity method, confidence is maintained in selecting the most important and correct content in an instrument, which is quantified by content validity ratio (CVR). The CVR is an item statistic that is useful in the rejection or retention of specific items. After items have been identified for inclusion in the final form, the content validity index (CVI) is computed for the whole test. In CVR, experts are requested to specify whether an item is necessary for operating a construct in a set of items or not. To this end, they are requested to score each item from 1 to 3 in with a three-degree range of "*not necessary, useful but not essential, essential*". CVR varies between 1 and −1. Greater levels of content validity exist as larger numbers of panelists agree that a particular item is essential. Using these assumptions, Lawshe developed a formula termed the CVR as: $CVR = (N_e − N/2)/(N/2)$, in which the $N_e$ is the number of panelists indicating "essential" and N is the total number of panelists. The numeric value of content validity ratio is determined by Lawshe table. In validating an instrument, then, a CVR value is computed for each item. From these are eliminated those items in which concurrence by members of Panel might reasonably have occurred through chance. Schipper (this table had been calculated for Lawshe by his friend, Lowell Schipper) has provided the data from which table 1 was prepared. Note, for example, that when a content evaluation panel is composed of 15 members, a minimum CVR of 0.49 is required to satisfy the 5% level. Only those items with CVR values meeting this minimum are retained in the final form of the instrument (27).

While the CVR is a direct linear transformation from the percentage saying "essential", its utility derives from its characteristics:

- When fewer than half say "essential", the CVR is negative.
- When half say "essential" and half do not, the CVR is 0.
- When all say "essential", the CVR is computed to be 1.00 (it is adjusted to 0.99 for ease of manipulation).
- When the number saying "essential" is more than half, but less than all, the CVR is

somewhere between 0 and 0.99 (27, 28).

In nursing the most widely reported approach for content validity is the CVI (3, 12, 25). A panel of content expert is asked to rate each instrument item in terms of clarity and its relevancy to the construct underlying study as per the theoretical definitions of the construct itself and its dimensions on a 4-point ordinal scale (1 = not relevant, 2 = somewhat relevant, 3 = quite relevant, 4 = highly relevant) (25). A table like the one shown below was added to the letter of request to guide experts for scoring method.

| Relevancy | Clarity |
|---|---|
| 1 = not relevant | 1 = not clear |
| 2 = item need some revision | 2 = item need some revision |
| 3 = relevant but need minor revision | 3 = clear but need minor revision |
| 4 = very relevant | 4 = very clear |

To obtain CVI for relevancy and clarity of each item (item levels [I-CVIs]), the number of those judging the item as relevant or clear (rating 3 or 4) was divided by the number of content experts, but for relevancy, CVI can be calculated both for I-CVIs and the scale-level (S-CVI). In item level, I-CVI is computed as the number of experts giving a rating 3 or 4 to the relevancy of each item, divided by the total number of experts. The I-CVI expresses the proportion of agreement on the relevancy of each item, which is between 0 and 1 (3, 29) and the S-CVI is defined as "the proportion of total items judged content valid" (3) or "the proportion of items on an instrument that achieved a rating of 3 or 4 by the content experts" (18).

Although instrument developers almost never give report what method have used to computing the scale-level index of an instrument (S-CVI) (6), there are two methods for calculating it, one method requires universal agreement (UA) among experts (S-CVI/UA), but a less conservative method is averages (Ave) the item-level CVIs (S-CVI/Ave). For calculating them, first, the scale is dichotomized by combining values 3 and 4 together and 2 and 1 together and two-choice options including "*relevant* and *not relevant*" are formed for each item (3, 25). Then, in the universal agreement approach, the number of items considered *relevant* by all the judges (or number of items with CVI equal to 1) is divided by the total number of items. In the average approach, the sum of I-CVIs is divided by the total number of items (10). Table 2 provides an example for better understanding about calculation CVI for the items of an instrument and S-CVI for the instrument by both methods. As the values obtained from both methods might be different, instrument makers should mention the method used for calculating it (6). Davis proposes that researchers should consider 80% agreement or higher among judges for new instruments (25). Judgment on each item is made as follows: If the I-CVI is higher than 79%, the item will be appropriate. If it is between 70% and 79%, it needs revision. If it is less than 70%, it is eliminated (30).

Although CVI is extensively used to estimate content validity by researchers, this index does not consider the possibility of inflated values because of the chance agreement. Therefore, CVI and the Kappa coefficient of agreement can provide quantifiable methods for evaluating the judgments of content experts. Kappa offers additional information beyond proportion agreement because it removes random chance agreement. For a better understanding of inter-rater agreement in general, and to increase confidence in the content validity of new instruments, researchers should report both the proportion agreement, as an indication of data variability, and the Kappa as a measure of agreement beyond chance (31). In other words, Kappa statistic is a consensus index of inter-rater agreement that adjusts for chance agreement (10) and is an important supplement to CVI because Kappa provides information about the degree of agreement beyond chance (7). Nevertheless, CVI is mostly used by researchers because it is simple for calculation, easy to understand and provide information about each item, which can be used for modification or deletion of instrument items (6, 10).

To calculate modified Kappa statistic, the probability of chance agreement was first calculated for each item by the following formula:

$$P_C = [N!/A!(N - A)!] * 0.5^N$$

**Table 1.** Minimum values of content validity ratio and one-tailed test, P = 0.05

| Number of panelists | Minimum value | Number of panelists | Minimum value |
|---|---|---|---|
| 5 | 0.99 | 11 | 0.59 |
| 6 | 0.99 | 12 | 0.56 |
| 7 | 0.99 | 13 | 0.54 |
| 8 | 0.75 | 14 | 0.51 |
| 9 | 0.78 | 15 | 0.49 |
| 10 | 0.62 | 20 | 0.42 |

**Table 2.** An example for calculating item-level content validity index and scale-level content validity index by two approaches of universal agreement of scale-level content validity index and averages of scale-level content validity index

| Items | (rating 3 or 4) | (rating 1 or 2) | I-CVIs | Interpretation |
|---|---|---|---|---|
| 1 | 14 | 0 | 1.000 | Appropriate |
| 2 | 12 | 2 | 0.857 | Appropriate |
| 3 | 13 | 1 | 0.928 | Appropriate |
| 4 | 12 | 2 | 0.857 | Appropriate |
| 5 | 11 | 3 | 0.785 | Need for revision |
| 6 | 14 | 0 | 1.000 | Appropriate |
| 7 | 12 | 2 | 0.857 | Appropriate |
| 8 | 8 | 6 | 0.571 | Eliminated |
| 9 | 14 | 0 | 1.000 | Appropriate |
| Number of items considered relevant by all the judges = 3 | | | | |
| Number of terms = 9 | | | S-CVI/Ave or average of I-CVIs = 0.872 | |
| S-CVI/UA = 3/9 = 0.333 | | | | |

Note: Number of experts = 14. Interpretation of I-CVIs: If the I-CVI is higher than 79%, the item will be appropriate. If it is between 70% and 79%, it needs revision. If it is less than 70%, it is eliminated. I-CVI: Item-level content validity index; S-CVI: Scale-level content validity index; UA: Universal agreement; Ave: Average

In this formula, N = number of experts in a panel and A = number of panelists who agree that the item is relevant.

After calculating I-CVI for all instrument items, finally, Kappa was computed by entering the numerical values of probability of chance agreement ($P_C$) and CVI of each item (I-CVI) in the following formula:

K = (I-CVI − $P_C$)/(1 − $P_C$).

Evaluation criteria for Kappa is the values above 0.74, between 0.60 and 0.74, and the ones between 0.40 and 0.59 are considered as excellent, good, and fair, respectively (32).

Polit et al. states that after controlling items by calculating modified Kappa statistic, each item with I-CVI equal or higher than 0.78 would be considered excellent. Researchers should note that, as the number of experts in panel increases, the probability of chance agreement diminishes and values of I-CVI and Kappa converge (10).

Requesting panel members to evaluate instrument in terms of comprehensiveness would be the last step of measuring the content validity. This step is necessary because an instrument may have acceptable inter-rater agreement, but still not cover the content domain. In judging the entire instrument, content experts evaluate whether the complete set of instrument items is sufficient to represent the total content domain. Is it needed to eliminate or add any item? According to members' judgment, proportion of agreement is calculated for the comprehensiveness of each dimension and the entire instrument. In order to the number of experts who have identified instrument comprehensiveness as favorable is divided into the total number of experts (3, 12).

***Determining Face Validity of an Instrument***

Face validity is used as a supplemental form of validity, supporting content validity, and answers this question whether an instrument apparently has validity for subjects, patients, and/or other participants. Face validity concerns judgments about items after an instrument is constructed, whereas content validity is more properly ensured by the plan of content and item construction before it is constructed. Thus, face validity can be considered as one limited aspect of content validity, concerning an inspection of the final product sure that nothing went wrong transforming plans into a completed instrument.

Face validity means if the designed instrument is apparently related to the construct underlying study. Do participants agree with items and wording of them in an instrument to realize research objectives? Face validity is related to the appearance and apparent attractiveness of an instrument, which may affect the instrument acceptability by respondents (11). It is obtained when the instrument users and subjects under study recognize that the instrument is suitable for measuring pertinent attributes. In principle, face validity is not considered as validity as far as measurement principles are concerned. In fact, it does not consider what to measure, but it focuses on the appearance of instrument (9). The overall validity of an instrument changes as words and items in face validity change. Therefore, determining face validity should be considered as the first measure (33).

To determine face validity of an instrument, researchers use respondents and experts' viewpoints. Difficulty level of items, desired suitability and relationship between items and the main objective of an instrument, ambiguity and misinterpretations of items, and/or incomprehensibility of the meaning of words are the issues discussed in the interviews (34).

Although experts play a vital role in content validity, instrument review by a sample of subjects drawn from the target population is another important component of content validation. These individuals are asked to review instrument items because of their familiarity with the construct through direct personal experience (12). Furthermore, they will be asked to identify the items they thought are the most important for them, and grade their importance on a 5-point Likert scale including very important (5), important (4), relatively important (3), slightly important (2), and unimportant. In quantities method, for calculation item impact score, the first is calculated percent of patients who scored 4 or 5 to item importance (frequency), and the mean importance score of item (importance) and then item impact score of instrument items was calculated by the following formula:

Item Impact Score = Frequency * Importance

If the item impact of an item is equal to or greater than 1.5 (which corresponds to a mean frequency of 50% and mean importance of 3 on the 5-point Likert scale), it is maintained in the instrument; otherwise, it is eliminated (35).

Finally, it should be said that validation is a lengthy process, in the first-step of which, the content validity should be studied and the following analyses should be directed include reliability evaluation (through internal consistency and test-retest), construct validity (through factor analysis) and criterion-related validity (12).

Some limitations of content validity studies should be noted. Experts' feedback is subjective; thus, the study is subjected to bias that may exist among the experts. If content domain is not well-identified, this type of study does not necessarily identify content that might have been omitted from the instrument. However, experts are asked to suggest other items for the instrument, which may help minimize this limitation (11).

## Conclusion

Content validity study is a systematic, subjective and two-stage process. In the first stage, instrument development is carried out and in the second stage, judgment/quantification on instrument items is performed, and content experts study the accordance between theoretical and operational definitions. Such process should be the leading study in the process of making instrument to guarantee instrument reliability and prepare a valid instrument in terms of content for preliminary test phase.

Understanding content validity is important for nursing researchers because they should realize if the instruments they use for their studies are suitable for the construct, population under study, and sociocultural background in which the study is carried out, or there is a need for new or modified instruments. Training on content validity study helps students, researchers, and clinical staffs better understand, use and criticize research instruments with a more accurate approach.

## References

1. McGartland Rubio D, Berg-Weger M, Tebb SS, Lee ES, Rauch S. Objectifying content validity: Conducting a content validity study

in social work research. Social Work Research 2003; 27(2): 94-104.

2. DeVon HA, Block ME, Moyle-Wright P, Ernst DM, Hayden SJ, Lazzara DJ, et al. A psychometric toolbox for testing validity and reliability. J Nurs Scholarsh 2007; 39(2): 155-64.

3. Lynn MR. Determination and quantification of content validity. Nurs Res 1986; 35(6): 382-5.

4. Waltz C, Strickland OL, Lenz E. Measurement in Nursing and Health Research. 4th ed. New York, NY: Springer Publishing Company; 2010. p. 163.

5. Newman I, Lim J. Content validity using a mixed methods approach: its application and development through the use of a table of specifications methodology. Journal of Mixed Methods Research 2013; 7(3): 243-60.

6. Polit DF, Beck CT. The content validity index: are you sure you know what's being reported? Critique and recommendations. Res Nurs Health 2006; 29(5): 489-97.

7. Wynd CA, Schmidt B, Schaefer MA. Two quantitative approaches for estimating content validity. West J Nurs Res 2003; 25(5): 508-18.

8. Yaghmale F. Content validity and its estimation. Journal of Medical Education 2003; 3(1): 25-7.

9. Anastasi A. Psychological testing. 6th ed. New York, NY: Macmillan; 1988.

10. Polit DF, Beck CT, Owen SV. Is the CVI an acceptable indicator of content validity? Appraisal and recommendations. Res Nurs Health 2007; 30(4): 459-67.

11. Nunnally JC, Bernstein IH. Psychometric theory. 3rd ed. Ventura, CA: Cram101 Incorporated, 2006.

12. Grant JS, Davis LL. Selection and use of content experts for instrument development. Res Nurs Health 1997; 20(3): 269-74.

13. Beck CT. Content validity exercises for nursing students. J Nurs Educ 1999; 38(3): 133-5.

14. Rattray J, Jones MC. Essential elements of questionnaire design and development. J Clin Nurs 2007; 16(2): 234-43.

15. Carmines EG, Zeller RA. Reliability and validity assessment. Thousand Oaks, CA: Sage Publications; 1979.

16. Cronbach LJ, Thornton GC. Test items to accompany Essentials of Psychological Testing. 3rd ed. New York, NY: Harper & Row; 1970.

17. Messick S. Evidence and Ethics in the Evaluation of Tests. Educational Researcher 1981; 10(9): 9-20.

18. Beck CT, Gable RK. Ensuring content validity: an illustration of the process. J Nurs Meas 2001; 9(2): 201-15.

19. Wilson HS. Research in nursing. Redwood City, CA: Addison-Wesley; 1989.

20. Tilden VP, Nelson CA, May BA. Use of qualitative methods to enhance content validity. Nurs Res 1990; 39(3): 172-5.

21. Ridenour CS, Benz CR, Newman I. Mixed methods research: exploring the interactive continuum. Carbondale, IL: SIU Press; 2008.

22. Priest J, McColl BA, Thomas L, Bond S. Developing and refining a new measurement tool. Nurse Researcher 1995; 2: 69-81.

23. Bowling A. Research methods in health: investigating health and health services. London, UK: Open University Press; 1997.

24. Grant JS, Kinney MR. Using the Delphi technique to examine the content validity of nursing diagnoses. Nurs Diagn 1992; 3(1): 12-22.

25. Lindsey Davis L. Instrument review: Getting the most from a panel of experts. Applied Nursing Research 1992; 5(4): 194-7.

26. Safikhani S, Sundaram M, Bao Y, Mulani P, Revicki DA. Qualitative assessment of the content validity of the Dermatology Life Quality Index in patients with moderate to severe psoriasis. J Dermatolog Treat 2013; 24(1): 50-9.

27. Lawshe CH. A quantitative approach to content validity. Personnel Psychology 1975; 28(4): 563-75.

28. Wilson FR. Recalculation of the critical values for Lawshe's content validity ratio. Measurement and Evaluation in Counseling and Development 2012; 45(3): 197-210.

29. Waltz CF, Bausell RB. Nursing research:

design, statistics, and computer analysis. Philadelphia, PA: F.A. Davis Co.; 1981.

30. Abdollahpour I, Nedjat S, Noroozian M, Majdzadeh R. Performing content validation process in development of questionnaires. Iran J Epidemiol 2010; 6(4): 66-74. [In Persian].

31. Brennan PF, Hays BJ. The kappa statistic for establishing interrater reliability in the secondary analysis of qualitative clinical data. Res Nurs Health 1992; 15(2): 153-8.

32. Cicchetti DV, Sparrow SA. Developing criteria for establishing interrater reliability of specific items: applications to assessment of adaptive behavior. Am J Ment Defic 1981; 86(2): 127-37.

33. Asadi-Lari M, Packham C, Gray D. Psychometric properties of a new health needs analysis tool designed for cardiac patients. Public Health 2005; 119(7): 590-8.

34. Banna JC, Vera Becerra LE, Kaiser LL, Townsend MS. Using qualitative methods to improve questionnaires for Spanish speakers: assessing face validity of a food behavior checklist. J Am Diet Assoc 2010; 110(1): 80-90.

35. Lacasse Y, Godbout C, Series F. Health-related quality of life in obstructive sleep apnoea. Eur Respir J 2002; 19(3): 499-503.